

AD-A096 813 NAVAL UNDERWATER SYSTEMS CENTER NEW LONDON CT NEW LO--ETC F/G 9/2
A FUNCTIONALLY EXPANDED COMPUTER PROGRAM FOR COMPUTING MACHINE---ETC(U)
FEB 81 M J GOLDSTEIN

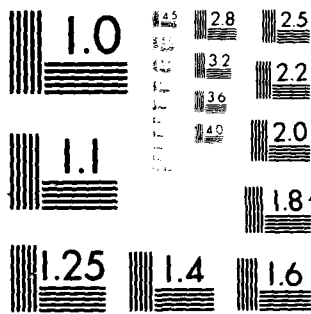
UNCLASSIFIED NUSC-TR-6421

NL

For I
S...



END
DATE
FILMED
4-5-81
DTIC



MICROCOPY RESOLUTION TEST CHART
NATIONAL BUREAU OF STANDARDS 1963-A

NUSC Technical Report 6421
12 February 1981

LEVEL II

12

A Functionally Expanded Computer Program for Computing Machine-Dependent Constants

M.J. Goldstein
Information Services Department

AD A 096813

DTIC
ELECTED
MAR 25 1981
S A



Naval Underwater Systems Center
Newport, Rhode Island / New London, Connecticut

Approved for public release;
distribution unlimited.

DTIC FILE COPY

81 3 25 013

Preface

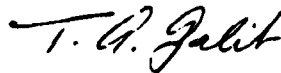
This report was prepared under NUSC Project No. 771Y00 for Special Projects and Studies.

The Technical Reviewer for this report was Richard Johnson (Code 711).

Acknowledgment

The author is grateful to James Ferrie (Code 325), Rosemary Molino (Code 325), and Nancy Sulinski (Code 7122) for running the computer program contained in Appendix B on the SEL 32/55, ITTEL AS/5, and the VAX 11/780. Furthermore, appreciation is due to John Lawson (Code 7122) for developing the algorithm for page length computation.

Reviewed and Approved: 12 February 1981



T. Galib
Information Services Department

The author of this report is located at the New London
Laboratory, Naval Underwater Systems Center,
New London, Connecticut 06320.

14

REPORT DOCUMENTATION PAGE		READ INSTRUCTIONS BEFORE COMPLETING FORM															
1. REPORT NUMBER TR-6421	2. GOVT ACCESSION NO. AD-A096813	3. RECIPIENT'S CATALOG NUMBER															
4. TITLE AND Subtitle A FUNCTIONALLY EXPANDED COMPUTER PROGRAM FOR COMPUTING MACHINE-DEPENDENT CONSTANTS.		5. TYPE OF REPORT & PERIOD COVERED															
7. AUTHOR(s) M. J. Goldstein		6. PERFORMING ORG. REPORT NUMBER															
9. PERFORMING ORGANIZATION NAME AND ADDRESS Naval Underwater Systems Center New London Laboratory New London, CT 06320		8. CONTRACT OR GRANT NUMBER(s) Technical rept.															
11. CONTROLLING OFFICE NAME AND ADDRESS Naval Underwater Systems Center Newport, RI 02840		10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS 771Y00															
14. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office) 12/51		12. REPORT DATE 12 February 1981															
		13. NUMBER OF PAGES 18															
		15. SECURITY CLASS. (of this report) UNCLASSIFIED															
		16. DECLASSIFICATION/DOWNGRADING SCHEDULE															
16. DISTRIBUTION STATEMENT (of this Report) Approved for public release; distribution unlimited.																	
17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)																	
18. SUPPLEMENTARY NOTES																	
19. KEY WORDS (Continue on reverse side if necessary and identify by block number) <table border="0"> <tr> <td>machine-dependent constants</td> <td>number of base digits</td> <td>smallest positive</td> </tr> <tr> <td>floating-point number</td> <td>number of bits in</td> <td>computer number</td> </tr> <tr> <td>unit roundoff error</td> <td>a digit</td> <td>largest computer number</td> </tr> <tr> <td>base</td> <td>effective bit precision</td> <td>program transportability</td> </tr> <tr> <td>mantissa</td> <td>page length</td> <td></td> </tr> </table>			machine-dependent constants	number of base digits	smallest positive	floating-point number	number of bits in	computer number	unit roundoff error	a digit	largest computer number	base	effective bit precision	program transportability	mantissa	page length	
machine-dependent constants	number of base digits	smallest positive															
floating-point number	number of bits in	computer number															
unit roundoff error	a digit	largest computer number															
base	effective bit precision	program transportability															
mantissa	page length																
20. ABSTRACT (Continue on reverse side if necessary and identify by block number) <p>In this report we present an updated version of an earlier FORTRAN program that computed machine-dependent constants. The earlier program has been expanded to compute (1) the base (radix) of the computer number system from the system's unit roundoff, (2) the number of digits in the fraction of a floating-point number, and (3) the number of bits in a digit. Furthermore, the program has been expanded to compute the page length of virtual memory machines like the VAX 11/780.</p>																	

405918

20. (Continued)

> An illustration is given of how the updated program can be used as an aid in transporting from one computer to another mathematical software that contains machine-dependent floating-point arithmetic parameters.

Table of Contents

	Page
Introduction	1
Theory and Algorithms	1
Preliminaries	1
Computation of the Base	3
Algorithm B	5
Largest Computer Number	5
Algorithm M	6
Page Length Computation	7
Algorithm PL	7
Program Description	7
An Application in Program Transportability	9
Conclusion	9
References	9
Appendix A: Theory and Algorithm for Unit Roundoff, Base B and t Unknown	11
Appendix B: Computer Program for Calculating Machine-Dependent Constants	13
Appendix C: Double Precision Program for Calculating Bessel Functions $J_N(x)$ and $I_N(x)$	15

A

A Functionally Expanded Computer Program For Computing Machine-Dependent Constants

Introduction

In a previous paper [1], a computer program was discussed that can be used (independently of any computer vendor claims) as a benchmark to reveal floating-point arithmetic characteristics of a binary device computer having a power of 2 base, $B (= 2^k, k \text{ a positive integer})$, and in transporting mathematical computer software that depend on the machine's relative accuracy (computer unit roundoff error) and/or arithmetic range from one computer to another. Furthermore, since the relative accuracy depends on the minimum (effective) number of significant bits in the mantissa (fraction) of a floating-point number and the roundoff property of the computer arithmetic, the program also computed the effective number of bits in the mantissa and whether the arithmetic rounds or chops. However, the computer arithmetic's base B and the number of base B digits in the mantissa of a floating-point number were not computed.

In this report, we present an updated program that computes the base B from the unit roundoff error. The new program also computes the number of base B digits in the mantissa of a floating-point number and the number of bits, namely k , in a base B digit. The updated program calculates the page length of virtual memory computer systems that organize data in virtual memory as the VAX 11/780 does. In order to enhance the program's transportability between computer systems, adjustments have been made to the program code and the algorithm for computing the largest floating-point base B computer number.

Theory and Algorithms

Preliminaries

If s is a non-zero, base B , floating-point, t -digit computer number, we assume its representation is

$$s = \pm B^e f, \quad -p \leq e \leq P$$

$$f = .b_1 b_2 \dots b_t, \quad b_i \in \{0, 1, 2, \dots, B-1\}, \quad b_1 \neq 0, \quad (1)$$

where B is a positive integer power of 2, namely 2^k , e is in a certain integer range with both p and P positive integers greater than t , and the value of f is given by

$$f = \sum_{i=1}^t b_i / B^i. \quad (2)$$

In other words, s is a real number that can be expressed as a normalized ($b_1 \neq 0$) base B number using at most t digits. In particular, if

$$b_i = \begin{cases} 1 & \text{if } i = 1, \\ 0 & \text{otherwise,} \end{cases} \quad (3)$$

then f assumes its smallest value:

$$f = B^{-1}. \quad (4)$$

If $b_i = B-1$ for $i = 1, 2, \dots, t$, then f assumes its largest value:

$$f = 1 - B^{-1}. \quad (5)$$

We assume that the internal computer representation of each base B digit b_i is given by k two state devices d_j , where each device takes on the value 0 or 1; that is,

$$\begin{aligned} b_i &= (d_k^{(i)} d_{k-1}^{(i)} \dots d_1^{(i)})_2 \\ &= \sum_{j=1}^k d_j^{(i)} 2^{j-1}, \quad d_j^{(i)} \in \{0, 1\}. \end{aligned} \quad (6)$$

In particular, if $f = 2^{-r}$ for $r = 1, 2, \dots, k$, then

$$\begin{aligned} b_1 &= 2^{k-r} = (\overbrace{0 \dots 0}^{r-1} 1 \overbrace{0 \dots 0}^{k-r})_2 \\ b_i &= 0 = (\overbrace{0 \dots 0}^k)_2, \quad i = 2, \dots, t. \end{aligned} \quad (7)$$

Therefore, the number of bits in the fraction (mantissa) f is kt . Thus, for a floating-point word of fixed length, say N bits with kt bits in the mantissa, $N-kt-1$ bits are reserved for the exponent e (stored as the excess quantity $e+a$, $a \geq 0$) and the remaining bit is reserved for the sign of f . Furthermore, since $b_1 \neq 0$, there can be as few as $kt-k+1$ significant (effective) bits in the fraction f .

By (1) and (4) the smallest positive real number in the computer number set is

$$m = B^{-p} B^{-1}, \quad (8)$$

whereas the largest positive real number is

$$M = B^p (1 - B^{-1}) \quad (9)$$

by (1) and (5). In a computation, the magnitude of results smaller than m produce underflow (usually returning a base B , t -digit representation for zero); the magnitude of results larger than M produce overflow.

The computer unit roundoff error u , which measures how closely a real number y can be approximated by a computer number s , is related to the error in s by the inequality

$$|y-s| \leq |s|u, \quad (10)$$

where

$$u = \begin{cases} B^{1-t} & \text{(chops)} \\ 0.5B^{1-t} & \text{(rounds)}. \end{cases} \quad (11)$$

The computer arithmetic chops if all of the mantissa's low order digits b_i beyond the t -th are simply discarded in the normalized base B , infinite digit representation of y ; rounds if b_t is incremented by 1 whenever $b_{t+1} \geq 2^{k-1}$, before discarding the low order digits.

Computation of the Base

In the earlier program, we computed the unit roundoff error when $B (= 2^k)$ and t are unknown. (The theory and algorithm for this is in Appendix A.) The following Theorem and Corollaries lead to an algorithm for computing the base B from the unit roundoff error.

Theorem. In the interval $[B^{e-l}, B^e]$, where $-p \leq e \leq P$, the floating-point numbers are uniformly spaced with spacing B^{e-t} .

Proof. Let $B^e f$ be a floating-point number that is in the interval $[B^{e-l}, B^e]$, such that

$$f = .b_1 \dots b_t, \quad B^{-l} \leq f < 1. \quad (12)$$

To obtain the next floating-point number in the interval, we add the base B digit one to the last digit b_t of f ; that is, the next floating-point, base B , t -digit number in $[B^{e-l}, B^e]$ is

$$B^e (f + .b'_1 \dots b'_t), \quad (13)$$

where

$$b'_i = \begin{cases} \overbrace{(0 \dots 0)}^k_2 & \text{if } i = 1, 2, \dots, t-1, \\ \overbrace{(0 \dots 0 1)}^{k-1}_2 & \text{if } i = t, \end{cases}$$

and

$$.b'_1 \dots b'_t = \sum_{i=1}^t b'_i (2^k)^i = (2^k)^{-t} = B^{-t}. \quad (14)$$

This completes the proof.

Corollary 1. The floating-point, base B , t -digit numbers in the interval $[B^{t-l}, B^t]$ are uniformly spaced with spacing $B^0 = 1$.

Proof. Immediate.

But the value of every floating-point, base B , t -digit number in $[B^{t-1}, B^t]$ is an integer. Therefore, every integer in $[B^{t-1}, B^t]$ can be represented exactly as a floating-point, base B , t -digit number, since B^{t-1} can be. In particular, the integers $2^L B^{t-1}$ for $L=0, 1, 2, \dots, k$, have exact floating-point, base B , t -digit representations.

Corollary 2. The floating-point, base B , t -digit numbers in the interval $[B^t, B^{t+1}]$ are uniformly spaced with spacing B .

Proof. Since $B^t = B^{(t+1)-1}$, set $e = t + 1$ in the Theorem. This completes the proof.

Thus, although B^t can be represented exactly as a floating-point, base B , t -digit number $B^t + 1$ cannot be. Then how is $B^t + 1$ approximated in the computer number set?

Consider that

$$\begin{aligned} B^t + 1 &= B^{t+1} B^{-1} + B B^{-1} \\ &= B^{t+1} B^{-1} + B^{t+1} (B^{-1} B^{-1}) \\ &= B^{t+1} (B^{-1} + B^{-t-1}) \\ &= B^{t+1} (. b_1 \dots b_t b_{t+1}), \end{aligned} \tag{15}$$

where

$$b_i = \begin{cases} 1 & \text{for } i=1, t+1 \\ 0 & \text{for } i=2, \dots, t. \end{cases}$$

Therefore, if the computer arithmetic rounds, the base B , t -digit, floating-point approximation to $B^t + 1$ is given by

$$B^t \oplus 1 = B^{t+1} . b_1 \dots b_{t-1} b_t^* \tag{16}$$

with

$$b_t^* = \begin{cases} 1 & \text{if } B=2, \\ 0 & \text{if } B=2^k \text{ for } k \geq 2. \end{cases}$$

On the other hand, if the computer arithmetic chops, then for $B=2^k$ ($k=1, 2, \dots$) the base B , t -digit, floating-point approximation to the sum $B^t + 1$ is given by (16) with $b_t^*=0$.

Hence, the computer arithmetic will approximate the real sum $B^t + 1$ by the computer number B^t , when the computer arithmetic rounds if $B=2^k$ for $k \geq 2$; if $B=2$, then the sum is approximated by $B^t + B$. On the other hand, if the computer arithmetic chops, it approximates the sum $B^t + 1$ by B^t for all values of B that are positive integer powers of two.

Therefore, given the unit roundoff error u in (11), since its reciprocal u^{-1} is an integer in the interval $[B^{t-1}, B^t]$, and

$$2^L u^{-1} = B^t$$

when

$$L = \begin{cases} k & \text{(chops)} \\ k-1 & \text{(rounds)}, \end{cases} \quad (17)$$

we have the following algorithm.

Algorithm B

Given the unit roundoff u , the effective precision, say NBIT, and whether the computer arithmetic rounds or chops, the algorithm for computing B , the binary length of a base B digit (namely k) and the number of base B digits in the mantissa (namely t) is:

B1. Set k to zero (written symbolically as $k \leftarrow 0$).

$$s \leftarrow u^{-1}$$

$$s_2 \leftarrow s$$

B2. Do while $(s_2 + 1 \neq s_2)$.

$$s_2 \leftarrow 2s_2$$

$$k \leftarrow k + 1$$

End Do while

B3. $B \leftarrow s_2/s$

B4. If (machine rounds and $B > 2$)

$$B \leftarrow 2B$$

$$k \leftarrow k + 1$$

End if

B5. $t \leftarrow (\text{NBIT} + k - 1)/k$

Largest Computer Number

In the earlier program [1] in order to compute the largest computer number

$$M = B^P (1 - B^{-t}), \quad (18)$$

we begin by taking the reciprocal of the smallest positive computer number

$$m = B^{-p-1}. \quad (19)$$

If the reciprocal does not cause overflow, the program outputs an adjusted value of the reciprocal as an approximation to M . Otherwise, we find the reciprocal of the product of m by the smallest power of two that does not cause overflow. This gives

$$M = B^P(0.5). \quad (20)$$

Then we output an adjusted value of this result as an approximation to M .

The problem with this algorithm is that it requires invoking a routine (not available on all systems) to test whether an overflow occurs. We have corrected this in the following way.

Assume that the relation between the smallest and largest exponents of B is given by

$$p = P + 1. \quad (21)$$

This is so on the UNIVAC 1108. On some machines the relation is $p = P$; e.g., VAX 11/780. Then with $p = P + 1$

$$m = B^{-p-1} = B^{-P-2} \quad (22)$$

and

$$m^{-1} = B^{P+2} = B^{P+3} B^{-1}, \quad (23)$$

so that

$$(B^3 m)^{-1} (1 - B^{1-t}) B = B^P (1 - B^{1-t}) \quad (24)$$

approximates M without computer arithmetic overflow, provided the order of operation is performed from left to right as indicated in (24). For if one were to compute $(B^3 m)^{-1} B$ when $p = P + 1$, the result B^P would cause computer overflow.

Similarly, if $p = P$, then

$$(B^3 m)^{-1} (1 - B^{1-t}) B = B^{P-1} (1 - B^{1-t}) \quad (25)$$

so that the approximation is off by B .

Algorithm M

Thus, given m , B and t , the algorithm for approximating M without computer overflow is

- M1. Compute B^M and store it in M.
- M2. Compute the reciprocal of M and store the result in M.
- M3. Compute and store the product $M(1-B^{1-t})$ in M.
- M4. Compute and store the product MB in M. Output M.

Page Length Computation

The idea behind the algorithm for determining page length is based on the fact that local data and program code are stored by the VAX 11/780 on separate pages on secondary computer memory, with program code following local data.* Therefore, if the number of computer words, e.g., I, required by local data occupies less than a page of secondary memory, then the first non-zero location following the Ith local data word is a program instruction.

Algorithm PL

Thus, dimensioning an integer array, e.g., IPGLN, as having length I, and assuming the number of local data items is I, the algorithm for the page length is

- PL1. $NMBDT \leftarrow I + 1$
- PL2. $J \leftarrow NMBDT$
- PL3. Do while (IPGLN(J) = 0)
 - $J \leftarrow J + 1$
 End Do while
- PL4. $J \leftarrow J - 1$

The final value of J is the page length.

Program Description

The updated single precision program that incorporates Algorithms B, M and PL is given in Appendix B. Algorithm u in Appendix A for computing the unit roundoff error, which is encoded in the earlier program [1], is implemented in the new program without any logical changes.

Algorithm PL for computing page length is implemented in the program for systems like the VAX 11/780 with page length as great as 2048 words, where the local data occupies no more than 17 single precision words for the single precision

*Local data refers to data only within the program module that defines them; e.g., data items in DATA type statements and intermediate results calculated and stored by the program module.

version of the program, and no more than 29 single precision words for the double precision version of the program.

Unlike the earlier program in [1], the new program is written in ANSI 66 FORTRAN [2] (with the exception of the PRINT statement that is ANS 77 [3]) to promote compilability on computers that have different FORTRAN compilers. Furthermore, the program is designed so that it can be converted easily to double precision — simply remove the comment flag C from column 1 of comment lines 5, 6, 7, 86 and 87 in the program, and convert lines 8, 84 and 85 to comment lines.

Since the program requires no input data, it is convenient to use in determining the values of the machine-dependent parameters discussed in this report.

The program in Appendix B has compiled and executed successfully on the UNIVAC 1108, VAX 11/780, SEL 32/55 and ITEL AS/5.* The output for the double precision version of the program on the UNIVAC 1108 and VAX 11/780 appear in tables 1 and 2 below, where NBIT is the effective bit precision of a floating-point computer number; BASE is the computer arithmetic's base B; NBDGT is the number of base B digits in the fraction of a floating-point computer number; LNGH is the number of bits in a base B digit; and RNCHOP is the arithmetic's rounding property — 0 if the arithmetic chops results, 1 if it rounds. The number of bits in the mantissa of a floating-point computer number is the product of NBDGT by LNGH.

Table 1. UNIVAC 1108 Double Precision

```
UNIT ROUNDOFF ERROR = .173472348-17      NBIT = 60      RNCHOP = 0.

BASE = 2.  NBDGT = 60  LNGH = 1

SMALLEST F. P. NUMBER = .278134232313-308
APPROXIMATE LARGEST F. P. NUMBER = .898846567431+308
PAGE LENGTH = 29
```

Table 2. VAX 11/780 Double Precision

```
UNIT ROUNDOFF ERROR = 0.138777878E-16      NBIT = 56      RNCHOP = 1.

BASE = 2.  NBDGT = 56  LNGH = 1

SMALLEST F. P. NUMBER = 0.293873587706D-38
APPROXIMATE LARGEST F. P. NUMBER = 0.850705917302D+38
PAGE LENGTH = 126
```

*The page length computation does not work on the ITEL AS/5, since the ITEL system does not organize local data in virtual memory as the VAX does.

An Application in Program Transportability

Consider transporting a computer program from one computer to another, whose execution depends on the values of machine-dependent floating-point arithmetic parameters in the program. In general, if the parameters are not reset to the new machine's values, the transported program will not execute properly. The program in Appendix B, or its double precision version, may be used to determine the new parameter values.

For example, consider the segment of a double precision program [4] in Appendix C for calculating Bessel functions $J_\nu(x)$ and $I_\nu(x)$ of real argument and integer order, which is available on NUSC's UNIVAC 1108. The program contains several machine-dependent parameters that are explained in the program's commentary and are set to the double precision values of the UNIVAC 1108 in the DATA statement at line 75 of the program. Trying to execute this program on the VAX will fail as a result of computer arithmetic overflow if the DATA statement values of the parameters have not been modified to reflect the VAX's double precision arithmetic. These parameter values can be changed to reflect the VAX's values with the aid of the program in Appendix B. First, run the double precision version of the program in Appendix B on the VAX. Then from its output (in table 2) and the explanation of machine-dependent constants in the Bessel function subroutine in Appendix C, one obtains the following DATA statement for the VAX:

```
DATA NSIG,NTEN,LARGEX,EXPARG/17,37,10000,87.D0/
```

where EXPARG is the natural logarithm of the approximation to the largest floating-point number in table 2, rounded to the nearest integer.

Conclusion

The program for computing machine-dependent constants can be used as a benchmark for revealing features of floating-point number systems on binary device computers, where the base of the number system is a positive integer power of two, the mantissa is a normalized fraction and the computer arithmetic has at least one guard digit. Furthermore, the program can be used as an aid in transporting mathematical software containing machine-dependent parameters from one computer to another. In fact, it is possible to facilitate the transportation of such software between many different computers by incorporating in the software program instructions that automatically compute the appropriate machine values for the parameters.

References

1. Marvin J. Goldstein and Richard Johnson, "A FORTRAN Program for Determining Computer Arithmetic Characteristics," NUSC TM No. 781074, April 1978.
2. *American Standard FORTRAN X3.9 1966*, Business Equipment Manufacturers Association, March, 1966.

TR 6421

3. *American National Standard Programming Language FORTRAN, ANSI X3.9 1977*, American National Standards Institute, New York, NY, 1978.
4. David J. Sookne, "Bessel Functions of Real Argument and Integer Order," *J. Res. Nat. Bur. Stand. (U.S.) 77B (Mat. Sci)*, Nos. 3 & 4, 124-132 (July/December 1973).

Appendix A

Theory and Algorithm for Unit Roundoff: Base B and t Unknown

By Theorem 1, the spacing of floating-point, base B, t-digit computer numbers in $[B^0, B]$ is B^{1-t} . Therefore, $B^0 + B^{1-t}$ is in the computer number set; furthermore, the real numbers

$$B^0 + vB^{1-t}, \quad v = 1, 2, \dots, B^{t-1}(B-1)$$

are in $[B^0, B]$ and the base B, computer number set; and, therefore, so are the numbers

$$B^0 + 2^w B^{1-t}, \quad w = 0, 1, 2, \dots, k(t-1)-1,$$

where

$$2^w B^{1-t} = 2^{-L}, \quad L = k(t-1)-w,$$

are negative powers of two that are in the computer number set, by (7).

However,

$$\begin{aligned} y &= B^0 + 2^{-1} B^{1-t} = BB^{-1} + B^{1-t}(2^{k-1}/B) \\ &= B(B^{-1} + 2^{k-1}/B^{t+1}) \\ &= B(.b_1 \dots b_t b_{t+1}). \end{aligned}$$

with

$$b_i = \begin{cases} \overbrace{(0 \dots 0 \ 1)}^{k-1}_2 = 1 & \text{if } i = 1, \\ \overbrace{(0 \dots 0)}^k_2 = 0 & \text{if } i = 2, \dots, t, \\ \overbrace{(1 \ 0 \dots 0)}^{k-1}_2 = 2^{k-1} & \text{if } i = t+1, \end{cases}$$

is a real number that is not in the computer number set although $2^{-1} B^{1-t}$ is. The real number y is approximated by a computer number s , namely

$$y \doteq s = \begin{cases} B(.b_1 \dots b_t) & \text{(chops)} \\ B(.b_1 \dots b_t') & \text{(rounds)} \end{cases}$$

with

$$b_{t+1} = \underbrace{(0 \dots 0)}_{k-1} 1)_2$$

provided the computer has a guard digit in its arithmetic hardware registers for b_{t+1} . Hence

$$s = \begin{cases} 1 & \text{(chopped arithmetic)} \\ 1 + B^{t+1} & \text{(rounded arithmetic)}. \end{cases}$$

Thus, the algorithm for computing the computer arithmetic's relative accuracy u , effective precision, and rounding property is

Algorithm u

- u1. Compute $1 \oplus \epsilon (\epsilon = 2^{-t})$ for $t = 1, 2, \dots$, until either $1 \oplus \epsilon = 1$ or $(1 \oplus \epsilon) \ominus \epsilon \neq 1$.
- u2. If $(1 \oplus \epsilon) \ominus \epsilon \neq 1$, then $u = \epsilon$ and the machine rounds; otherwise $u = 2\epsilon$ and the machine chops.
- u3. The minimum number of significant bits in the mantissa is the last value of t computed in u1, namely $k(t-1) + 1$.

Appendix B

Computer Program for Calculating Machine-Dependent Constants

```

1  *****
2  C  READING F. F. WITH FEATURES BY M. J. GOLDSTEIN
3  C  *****
4  C  DIMENSION I(200)
5  C  DO WHILE (F(1) .GT. 0)
6  C  I(1) = F(1)
7  C  I(2) = F(2)
8  C  I(3) = F(3)
9  C  I(4) = F(4)
10 C  I(5) = F(5)
11 C  I(6) = F(6)
12 C  I(7) = F(7)
13 C  I(8) = F(8)
14 C  I(9) = F(9)
15 C  I(10) = F(10)
16 C  I(11) = F(11)
17 C  I(12) = F(12)
18 C  I(13) = F(13)
19 C  I(14) = F(14)
20 C  I(15) = F(15)
21 C  I(16) = F(16)
22 C  I(17) = F(17)
23 C  I(18) = F(18)
24 C  I(19) = F(19)
25 C  I(20) = F(20)
26 C  I(21) = F(21)
27 C  I(22) = F(22)
28 C  I(23) = F(23)
29 C  I(24) = F(24)
30 C  I(25) = F(25)
31 C  I(26) = F(26)
32 C  I(27) = F(27)
33 C  I(28) = F(28)
34 C  I(29) = F(29)
35 C  I(30) = F(30)
36 C  I(31) = F(31)
37 C  I(32) = F(32)
38 C  I(33) = F(33)
39 C  I(34) = F(34)
40 C  I(35) = F(35)
41 C  I(36) = F(36)
42 C  I(37) = F(37)
43 C  I(38) = F(38)
44 C  I(39) = F(39)
45 C  I(40) = F(40)
46 C  I(41) = F(41)
47 C  I(42) = F(42)
48 C  I(43) = F(43)
49 C  I(44) = F(44)
50 C  I(45) = F(45)
51 C  I(46) = F(46)
52 C  I(47) = F(47)
53 C  I(48) = F(48)
54 C  I(49) = F(49)
55 C  I(50) = F(50)
56 C  I(51) = F(51)
57 C  I(52) = F(52)
58 C  I(53) = F(53)
59 C  I(54) = F(54)
60 C  I(55) = F(55)
61 C  I(56) = F(56)
62 C  I(57) = F(57)
63 C  I(58) = F(58)
64 C  I(59) = F(59)
65 C  I(60) = F(60)
66 C  I(61) = F(61)
67 C  I(62) = F(62)
68 C  I(63) = F(63)
69 C  I(64) = F(64)
70 C  I(65) = F(65)
71 C  I(66) = F(66)
72 C  I(67) = F(67)
73 C  I(68) = F(68)
74 C  I(69) = F(69)
75 C  I(70) = F(70)
76 C  I(71) = F(71)
77 C  I(72) = F(72)
78 C  I(73) = F(73)
79 C  I(74) = F(74)
80 C  I(75) = F(75)
81 C  I(76) = F(76)
82 C  I(77) = F(77)
83 C  I(78) = F(78)
84 C  I(79) = F(79)
85 C  I(80) = F(80)
86 C  I(81) = F(81)
87 C  I(82) = F(82)
88 C  I(83) = F(83)
89 C  I(84) = F(84)
90 C  I(85) = F(85)
91 C  I(86) = F(86)
92 C  I(87) = F(87)
93 C  I(88) = F(88)
94 C  I(89) = F(89)
95 C  I(90) = F(90)
96 C  I(91) = F(91)
97 C  I(92) = F(92)
98 C  I(93) = F(93)
99 C  I(94) = F(94)
100 C  I(95) = F(95)
101 C  I(96) = F(96)
102 C  I(97) = F(97)
103 C  I(98) = F(98)
104 C  I(99) = F(99)
105 C  I(100) = F(100)
106 C  I(101) = F(101)
107 C  I(102) = F(102)
108 C  I(103) = F(103)
109 C  I(104) = F(104)
110 C  I(105) = F(105)
111 C  I(106) = F(106)
112 C  I(107) = F(107)
113 C  I(108) = F(108)
114 C  I(109) = F(109)
115 C  I(110) = F(110)
116 C  I(111) = F(111)
117 C  I(112) = F(112)
118 C  I(113) = F(113)
119 C  I(114) = F(114)
120 C  I(115) = F(115)
121 C  I(116) = F(116)
122 C  I(117) = F(117)
123 C  I(118) = F(118)
124 C  I(119) = F(119)
125 C  I(120) = F(120)
126 C  I(121) = F(121)
127 C  I(122) = F(122)
128 C  I(123) = F(123)
129 C  I(124) = F(124)
130 C  I(125) = F(125)
131 C  I(126) = F(126)
132 C  I(127) = F(127)
133 C  I(128) = F(128)
134 C  I(129) = F(129)
135 C  I(130) = F(130)
136 C  I(131) = F(131)
137 C  I(132) = F(132)
138 C  I(133) = F(133)
139 C  I(134) = F(134)
140 C  I(135) = F(135)
141 C  I(136) = F(136)
142 C  I(137) = F(137)
143 C  I(138) = F(138)
144 C  I(139) = F(139)
145 C  I(140) = F(140)
146 C  I(141) = F(141)
147 C  I(142) = F(142)
148 C  I(143) = F(143)
149 C  I(144) = F(144)
150 C  I(145) = F(145)
151 C  I(146) = F(146)
152 C  I(147) = F(147)
153 C  I(148) = F(148)
154 C  I(149) = F(149)
155 C  I(150) = F(150)
156 C  I(151) = F(151)
157 C  I(152) = F(152)
158 C  I(153) = F(153)
159 C  I(154) = F(154)
160 C  I(155) = F(155)
161 C  I(156) = F(156)
162 C  I(157) = F(157)
163 C  I(158) = F(158)
164 C  I(159) = F(159)
165 C  I(160) = F(160)
166 C  I(161) = F(161)
167 C  I(162) = F(162)
168 C  I(163) = F(163)
169 C  I(164) = F(164)
170 C  I(165) = F(165)
171 C  I(166) = F(166)
172 C  I(167) = F(167)
173 C  I(168) = F(168)
174 C  I(169) = F(169)
175 C  I(170) = F(170)
176 C  I(171) = F(171)
177 C  I(172) = F(172)
178 C  I(173) = F(173)
179 C  I(174) = F(174)
180 C  I(175) = F(175)
181 C  I(176) = F(176)
182 C  I(177) = F(177)
183 C  I(178) = F(178)
184 C  I(179) = F(179)
185 C  I(180) = F(180)
186 C  I(181) = F(181)
187 C  I(182) = F(182)
188 C  I(183) = F(183)
189 C  I(184) = F(184)
190 C  I(185) = F(185)
191 C  I(186) = F(186)
192 C  I(187) = F(187)
193 C  I(188) = F(188)
194 C  I(189) = F(189)
195 C  I(190) = F(190)
196 C  I(191) = F(191)
197 C  I(192) = F(192)
198 C  I(193) = F(193)
199 C  I(194) = F(194)
200 C  I(195) = F(195)
201 C  I(196) = F(196)
202 C  I(197) = F(197)
203 C  I(198) = F(198)
204 C  I(199) = F(199)
205 C  I(200) = F(200)
206 C  I(201) = F(201)
207 C  I(202) = F(202)
208 C  I(203) = F(203)
209 C  I(204) = F(204)
210 C  I(205) = F(205)
211 C  I(206) = F(206)
212 C  I(207) = F(207)
213 C  I(208) = F(208)
214 C  I(209) = F(209)
215 C  I(210) = F(210)
216 C  I(211) = F(211)
217 C  I(212) = F(212)
218 C  I(213) = F(213)
219 C  I(214) = F(214)
220 C  I(215) = F(215)
221 C  I(216) = F(216)
222 C  I(217) = F(217)
223 C  I(218) = F(218)
224 C  I(219) = F(219)
225 C  I(220) = F(220)
226 C  I(221) = F(221)
227 C  I(222) = F(222)
228 C  I(223) = F(223)
229 C  I(224) = F(224)
230 C  I(225) = F(225)
231 C  I(226) = F(226)
232 C  I(227) = F(227)
233 C  I(228) = F(228)
234 C  I(22
```

TR 6421

```
66      C
67      30 UNIT = UNIT*BASE**3
68      UNIT = ONE/UNIT
69      UNIT = UNIT*(ONE - BASE**(1 - NBDGT))
70      UNIT = UNIT*BASE
71      C CONTINUE
72      PRINT 60, UNIT
73      C
74      C ***DETERMINE PAGE LENGTH***
75      C
76      DO 35 J = NBDGT, 2050
77      IF (IPGLN(J).NE.0) GO TO 40
78      35 CONTINUE
79      40 J = J - 1
80      IF (J.LT.2049) PRINT 70, J
81      50 FORMAT(1H ,22HUNIT ROUNDOFF ERROR = ,E17.9,5X,7HNRBIT = ,I3,
82      X 5X,9HNRCHOP = ,F2.0 //)
83      52 FORMAT(1H ,7HBASE = ,F4.0,10H NBDGT = ,I4,9H LNGB = ,I4//)
84      55 FORMAT(1H ,24HSMALLEST F. P. NUMBER = ,E17.9)
85      60 FORMAT(1H ,35HAPPROXIMATE LARGEST F. P. NUMBER = ,E17.9)
86      C 55 FORMAT(1H ,24HSMALLEST F. P. NUMBER = ,D20.12)
87      C 60 FORMAT(1H ,35HAPPROXIMATE LARGEST F. P. NUMBER = ,D20.12)
88      70 FORMAT(1H ,14HPAGE LENGTH = ,I4)
89      STOP
90      END
```

14

Appendix C

Double Precision Program for Calculating Bessel Function $J_k(x)$ and $I_k(x)$

```

MULTIPLYING BY RESERVE
1      SUBROUTINE BESSEL(X, NR, IZE, R, NCALC)
2      C THIS ROUTINE CALCULATES BESSEL FUNCTIONS I AND J OF REAL
3      C ARGUMENT AND INTEGER ORDER.
4
5      C
6      C EXPLANATION OF VARIABLES IN THE CALLING SEQUENCE
7
8      C X DOUBLE PRECISION REAL ARGUMENT FOR WHICH I'S OR J'S
9      C ARE TO BE CALCULATED. IF I'S ARE TO BE CALCULATED,
10     C ABS(X) MUST BE LESS THAN EXFARG (WHICH SEE BELOW).
11     C NR INTEGER TYPE, 1 + HIGHEST ORDER TO BE CALCULATED.
12     C IT MUST BE POSITIVE.
13     C IZE INTEGER TYPE, ZERO IF J'S ARE TO BE CALCULATED, 1
14     C IF I'S ARE TO BE CALCULATED.
15     C R DOUBLE PRECISION VECTOR OF LENGTH NR, NEED NOT BE
16     C INITIALIZED BY USER. IF THE ROUTINE TERMINATES
17     C NORMALLY (NCALC = NR), IT RETURNS J(OR I)-SUB-ZERO
18     C THROUGH J(OR I)-SUB-NR-MINUS-ONE OF X IN THIS
19     C VECTOR.
20     C NCALC INTEGER TYPE, NEED NOT BE INITIALIZED BY USER.
21     C BEFORE USING THE RESULTS, THE USER SHOULD CHECK THAT
22     C NCALC=NR, I.E. ALL ORDERS HAVE BEEN CALCULATED TO
23     C THE DESIRED ACCURACY. SEE ERROR RETURNS BELOW.
24
25     C
26     C EXPLANATION OF MACHINE DEPENDENT CONSTANTS
27
28     C NSIG DECIMAL SIGNIFICANCE DESIRED. SHOULD BE SET TO
29     C IFIX(ALOG10(2)*NRBIT + 1), WHERE NRBIT IS THE NUMBER OF
30     C BITS IN THE MANTISSA OF A DOUBLE PRECISION VARIABLE.
31     C SETTING NSIG LOWER WILL RESULT IN DECREASED ACCURACY
32     C WHILE SETTING NSIG HIGHER WILL INCREASE CPU TIME
33     C WITHOUT INCREASING ACCURACY. THE TRUNCATION ERROR
34     C IS LIMITED TO  $7.5 \times 10^{**NSIG}$  FOR J'S OF ORDER LESS
35     C THAN ARGUMENT, AND TO A RELATIVE ERROR OF 1 FOR I'S
36     C AND THE OTHER J'S.
37     C NTEN LARGEST INTEGER N SUCH THAT  $10^{**N}$  IS MACHINE-
38     C REPRESENTABLE IN DOUBLE PRECISION.
39     C LARGEX UPPER LIMIT ON THE MAGNITUDE OF X. BEAR IN MIND
40     C THAT IF  $ABS(X) > N$ , THEN AT LEAST N ITERATIONS OF THE
41     C BACKWARD RECURSION WILL BE EXECUTED.
42     C EXFARG LARGEST DOUBLE PRECISION ARGUMENT THAT THE LIBRARY
43     C DEFP ROUTINE CAN HANDLE.
44
45     C
46     C ERROR RETURNS
47
48     C
49     C LET G DENOTE EITHER I OR J.
50     C IN CASE OF AN ERROR, NCALC.NE.NR, AND NOT ALL G'S
51     C ARE CALCULATED TO THE DESIRED ACCURACY.
52     C IF NCALC.LT.0, AN ARGUMENT IS OUT OF RANGE. NR.LE.0
53     C OR IZE IS NEITHER 0 OR 1 OR IZE-1 AND  $ABS(X) > EXFARG$ .
54     C IN THIS CASE, THE R VECTOR IS NOT CALCULATED, AND NCALC
55     C IS SET TO MIN0(NR+0)-1 SO NCALC.NE.NR.
56     C NR.GT.NCALC.GT.0 WILL OCCUR IF NR.GT.MAGX AND  $ABS(X)$ 
57     C  $SUB-NR-OF-X > SUB-MAGX-OF-X$ . I.E.  $10^{**NTEN}/2$ , I.E. NR
58     C IS MUCH GREATER THAN MAGX. IN THIS CASE, R(N) IS CALCU-
59     C LATED TO THE DESIRED ACCURACY FOR N.IE.NCALC, BUT FOR
60     C  $N>NCALC$ , N.IE.NR, PRECISION IS LOST. IF  $N.GT.NCALC$  AND
61     C  $ABS(R(NCALC)/R(N)) < 10^{**NSIG}$ , THEN ONLY THE FIRST NSIG-N
62     C SIGNIFICANT FIGURES OF R(N) MAY BE TRUSTED. IF THE USER
63     C WISHED TO CALCULATE R(N) TO HIGHER ACCURACY, HE SHOULD USE
64     C AN ASYMPTOTIC FORMULA FOR HIGH ORDER.
65

```

TR 6421

65 L DOUBLE PRECISION
66 1 X, B, P, TEST, TEMPA, TEMPB, TEMPC, EXPARG, SIGN, SUM, TOVER,
67 2 PLAST, FOLD, FSAVE, FSAVE1, ZERO, TENTH, HALF, QUART, ONE, TWO, TEN
68 3, SROOT, IUM
69
70 C 1, UNIT, RNDOP, SMALL, BIG, CMLOG
71 DIMENSION R(NB)
72 DATA ZERO, TENTH, QUART, HALF, ONE, TWO, TEN/0.0D0, .1D0, .25D0, .5D0,
73 1.0D0, 2.0D0, 10.0D0/
74 DATA NSIG, NTEN, LARGE, EXPARG/19, 307, 100000, 7.0D2/
75

THIS PAGE TO BE REMOVED
FROM COPY OF THIS MANUAL

Initial Distribution List

ADDRESSEE	NO. OF COPIES
Dep. USDR&E (Res. & Adv. Tech.) (R. M. Davis)	1
Dep. USDR&E (Dir. Elect. & Phys. Sc.) (L. Wiseberg)	1
OASN, Dept. Assit. Sec. (Res. & Adv. Tech.) (Dr. R. Hoglund)	1
ONR (ONR-200, -102, -212, -222, -230)	5
CNO (OP-02, -03-EG, -090, -098, -224, -35, -902, -951, -96, -961, -981G1, -981H, -982F, -983	14
CNM (MAT-08T2, -08T21, -08T24, SP-20, ASW-122, ASW-124)	6
DIA (DT-2C)	1
NAVSURFWEAPCTR, White Oak	1
DWTNSRDC ANNA	1
DWTNSRDC CARD	1
DTNSRDC BETHESDA	1
NRL	1
NRL, USRD	1
NRL, AESD	1
NORDA (Code 110, Dr. R. Goodman)	1
USOC (Code 241, 240)	2
NAVSUBASE, New London	1
NABSUBSUPACNLON	1
NOO (Code 02)	1
NAVELECSYSCOM (ELEX 00, 03, 304, 310, 5101, PME-108, -117)	7
NAVSEASYSYSCOM (SEA-003, -05R, -511, -06, -06D, -06R, -06V, -06Z, -61, -61R, -62, -62R, -63, -63R, -63R-1, -63R-13, -63X, -631X, -631Y, -632X, -902)	21
NAVSEASYSDET, Norfolk	1
NASC (AIR 610)	1
NAVAIRDEVCM (Code 00, 2052)	2
NOSC (Code 00, Code 6565)	2
NAVWPNSCEN	1
NAVCOASTSYSLAB	1
CIVENGRLAB	1
NAVSURFWPNCEN	1
CHESNAVFACENGCOM (FPO-IP3)	1
APL/UW, Seattle	1
ARL/Penn State	1
CENTER FOR NAVAL ANALYSES	1
DTIC	12
DARPA	1
NOAA/ERL	1
NATIONAL RESEARCH COUNCIL	1
WEAPON SYS EVAL GROUP	1
WOODS HOLE OCEANOGRAPHIC INST.	1
ENGINEERING SOCIETIES LIBRARY	1
MARINE PHYSICAL LAB, SCRIPPS	1

